

# **CAPSTONE PROJECT**

## **EDA ON HOTEL BOOKING ANALYSIS**

BY

CHETAN BHANGARE

NIKHIL BHANGARE

# PROBLEM STATEMENT

1. For this project we will be analyzing Hotel Booking data. This data set contains booking information for a city hotel and a resort hotel, and includes information such as when the booking was made, length of stay, the number of adults, children, and/or babies, and the number of available parking spaces.
2. Hotel industry is a very volatile industry and the bookings depends on above factors and many more.
3. The main objective behind this project is to explore and analyze data to discover important factors that govern the bookings and give insights to hotel management ,which can perform various campaigns to boost the business and performance.

- We will divide our work flow into three steps

Data Collection and  
Understanding

Data Cleaning and  
Manipulation

Exploratory Data  
Analysis(EDA)

- EDA will be divided into following 3 UBM rule analysis.
- 1) Univariate analysis: Univariate analysis is the simplest of the three analyses where the data you are analyzing is only one variable.
- 2) Bivariate analysis: Bivariate analysis is where you are comparing two variables to study their relationships.
- 3) Multivariate analysis: Multivariate analysis is similar to Bivariate analysis but you are comparing more than two variables

# Data Collection and Understanding

- □ After collecting data it's very important to understand your data. So we had hotel Booking analysis data.
- Which had 119390 rows and 32 columns. So let's understand this 32 columns.
- **Data Description:**
  - hotel :Resort Hotel or City Hotel
  - is\_canceled : Value indicating if the booking was canceled (1) or not (0)
  - lead\_time : Number of days that elapsed between the entering date of the booking and the arrival date

# Data Collection and Understanding

- arrival\_date\_year : Year of arrival date
- arrival\_date\_month : Month of arrival date
- arrival\_date\_week\_number : Week number of year for arrival date
- arrival\_date\_day\_of\_month : Day of arrival date
- stays\_in\_weekend\_nights : Number of weekend nights
- stays\_in\_week\_nights : Number of week nights.
- adults : Number of adults
- children : Number of children
- babies : Number of babies
- meal : Type of meal booked.
- country : Country of origin.

# Data Collection and Understanding

- market\_segment : Market segment designation. (TA/TO)
- distribution\_channel : Booking distribution channel.(T/A/TO)
- is\_repeated\_guest : is a repeated guest (1) or not (0)
- previous\_cancellations : Number of previous bookings that were cancelled by the customer prior to the current booking
- previous\_bookings\_not\_canceled : Number of previous bookings not cancelled by the customer prior to the current booking
- reserved\_room\_type : Code of room type reserved.
- assigned\_room\_type : Code for the type of room assigned to the booking.
- booking\_changes : Number of changes made to the booking from the moment the booking was entered on the PMS until the moment of check-in or cancellation

# Data Collection and Understanding

- deposit\_type : No Deposit, Non Refund , Refundable.
- agent : ID of the travel agency that made the booking
- company : ID of the company/entity that made the booking .
- days\_in\_waiting\_list : Number of days the booking was in the waiting list before it was confirmed to the customer
- customer\_type : type of customer. Contract,Group,transient,Transient party.
- adr : Average Daily Rate as defined by dividing the sum of all lodging transactions by the total number of staying nights
- required\_car\_parking\_spaces : Number of car parking spaces required by the customer
- total\_of\_special\_requests : Number of special requests made by the customer (e.g. twin bed or high floor)
- reservation\_status : Reservation last status

# Data Cleaning

```
duplicate_rows_df = df[df.duplicated()].shape  
  
print(f"the no. of duplicate rows : " , duplicate_rows_df)
```

the no. of duplicate rows : (31994, 32)

Lets drop the duplicate values

```
df=df.drop_duplicates()  
df.shape
```

(87396, 32)



# Data Cleaning

```
df.isnull().sum()

hotel                0
is_canceled          0
lead_time            0
arrival_date_year    0
arrival_date_month   0
arrival_date_week_number 0
arrival_date_day_of_month 0
stays_in_weekend_nights 0
stays_in_week_nights 0
adults               0
children             4
babies               0
meal                 0
country              452
market_segment       0
distribution_channel 0
is_repeated_guest    0
previous_cancellations 0
previous_bookings_not_canceled 0
reserved_room_type   0
assigned_room_type    0
booking_changes       0
deposit_type         0
agent                12193
company              82137
days_in_waiting_list 0
customer_type         0
adr                  0
required_car_parking_spaces 0
total_of_special_requests 0
reservation_status    0
reservation_status_date 0
dtype: int64
```

Since the column named **Company and Agents** have lots of null values , we will drop these columns

# Data Cleaning

```
df = df.drop(columns=['company','agent'])
```

```
df.isnull().sum()
```

```
hotel                0
is_canceled          0
lead_time            0
arrival_date_year    0
arrival_date_month   0
arrival_date_week_number 0
arrival_date_day_of_month 0
stays_in_weekend_nights 0
stays_in_week_nights 0
adults               0
children             4
babies               0
meal                 0
country              452
market_segment       0
distribution_channel 0
is_repeated_guest    0
previous_cancellations 0
previous_bookings_not_canceled 0
reserved_room_type   0
assigned_room_type    0
booking_changes       0
deposit_type          0
days_in_waiting_list 0
customer_type         0
adr                  0
required_car_parking_spaces 0
total_of_special_requests 0
reservation_status    0
reservation_status_date 0
dtype: int64
```

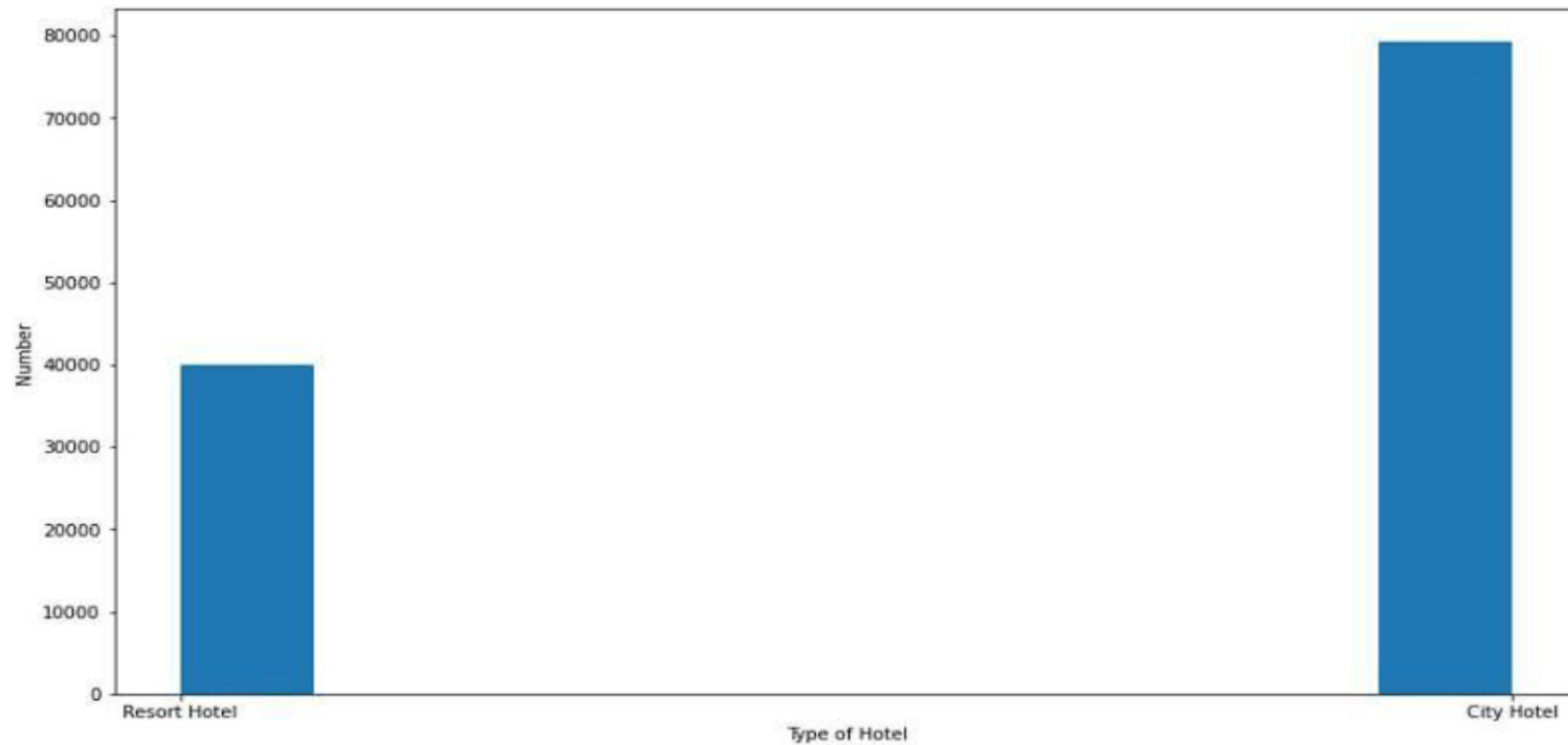
```
df=df.dropna()
```

```
#To insure we don't have any null values  
df.isnull().sum()
```

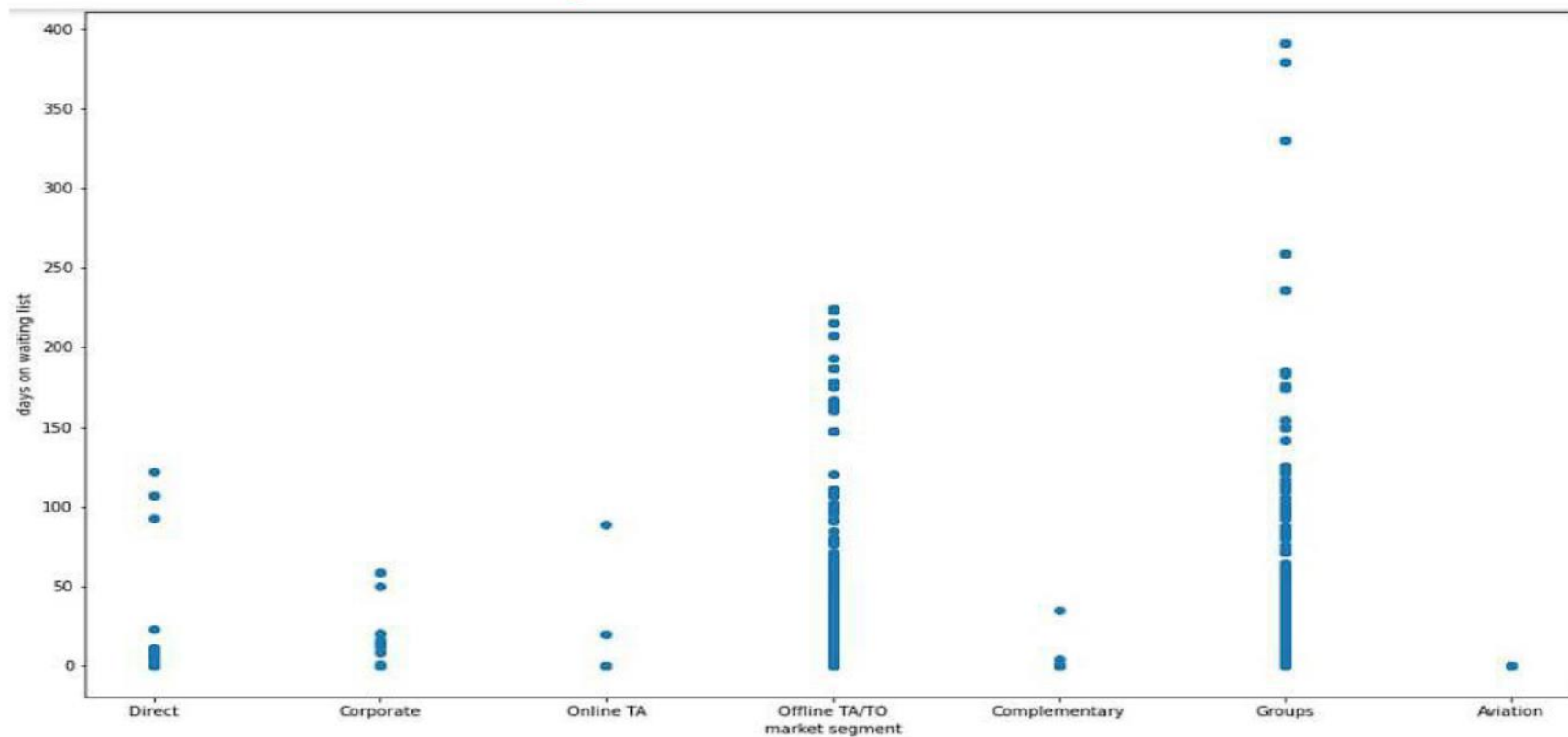
Since we only have children and country data with null values , so drop the unavailable data

# UNIVARIATE ANALYSIS

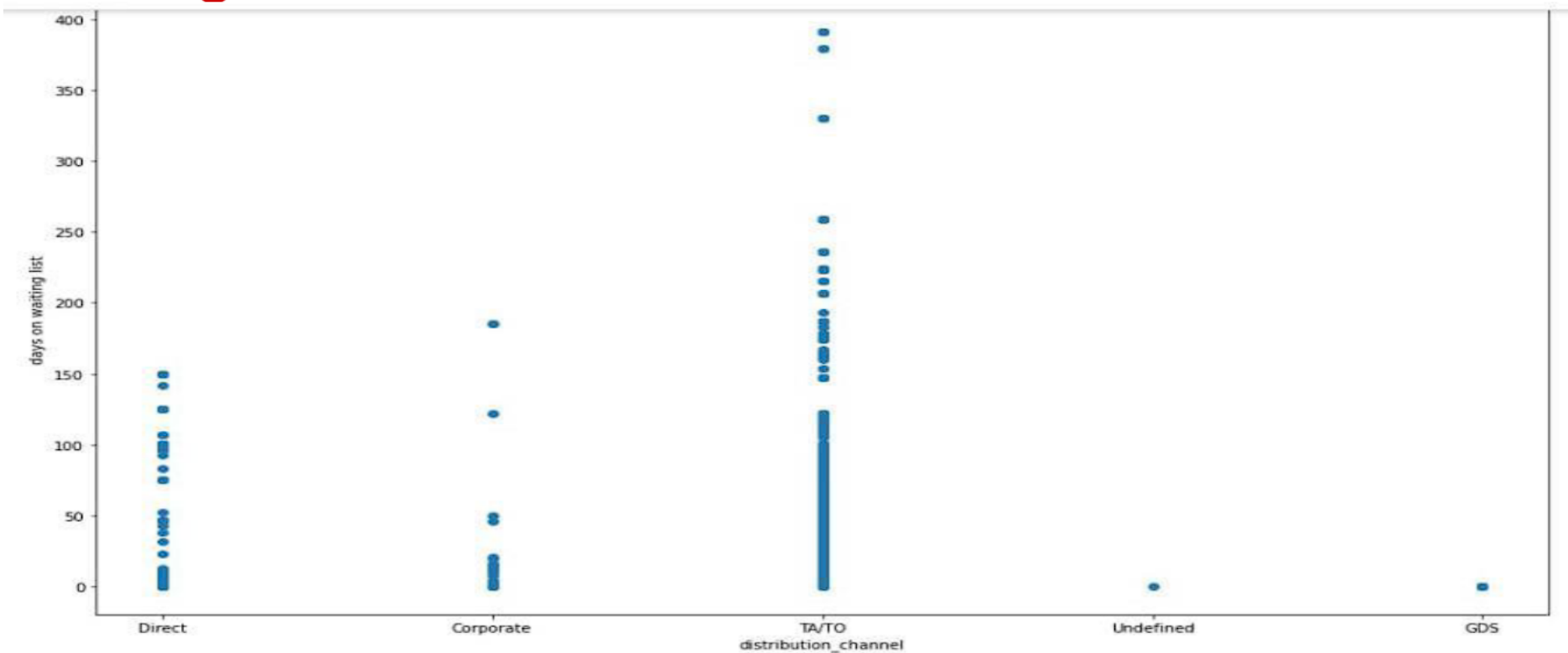
Number of Bookings for various types of hotels



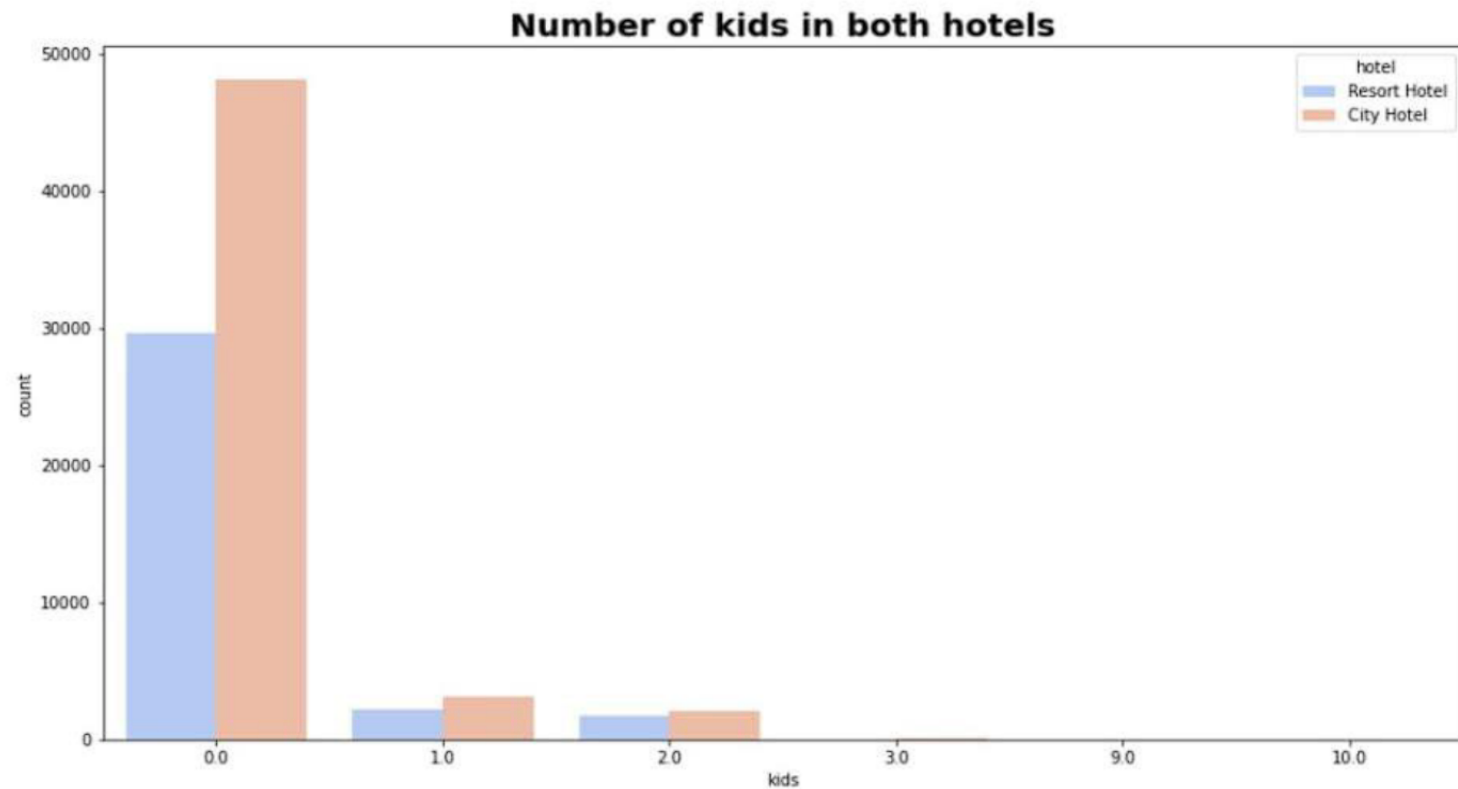
## Plot between Type of market segment and Waiting list for the booking



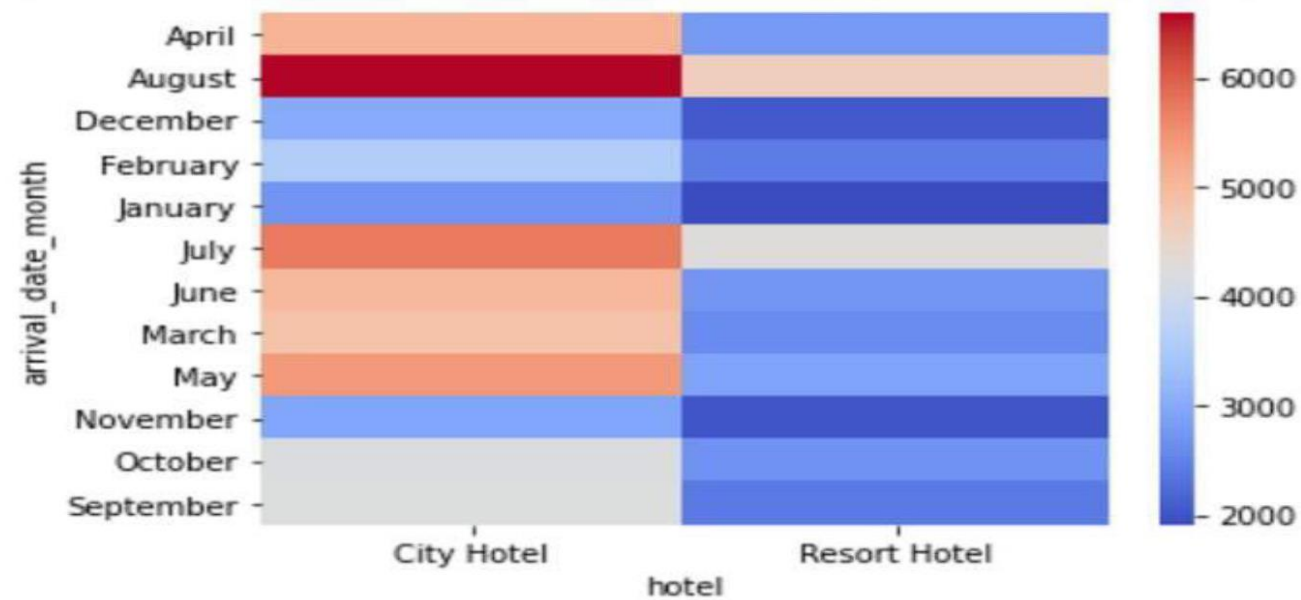
## Plot between Distributing Channel and Days on the waiting



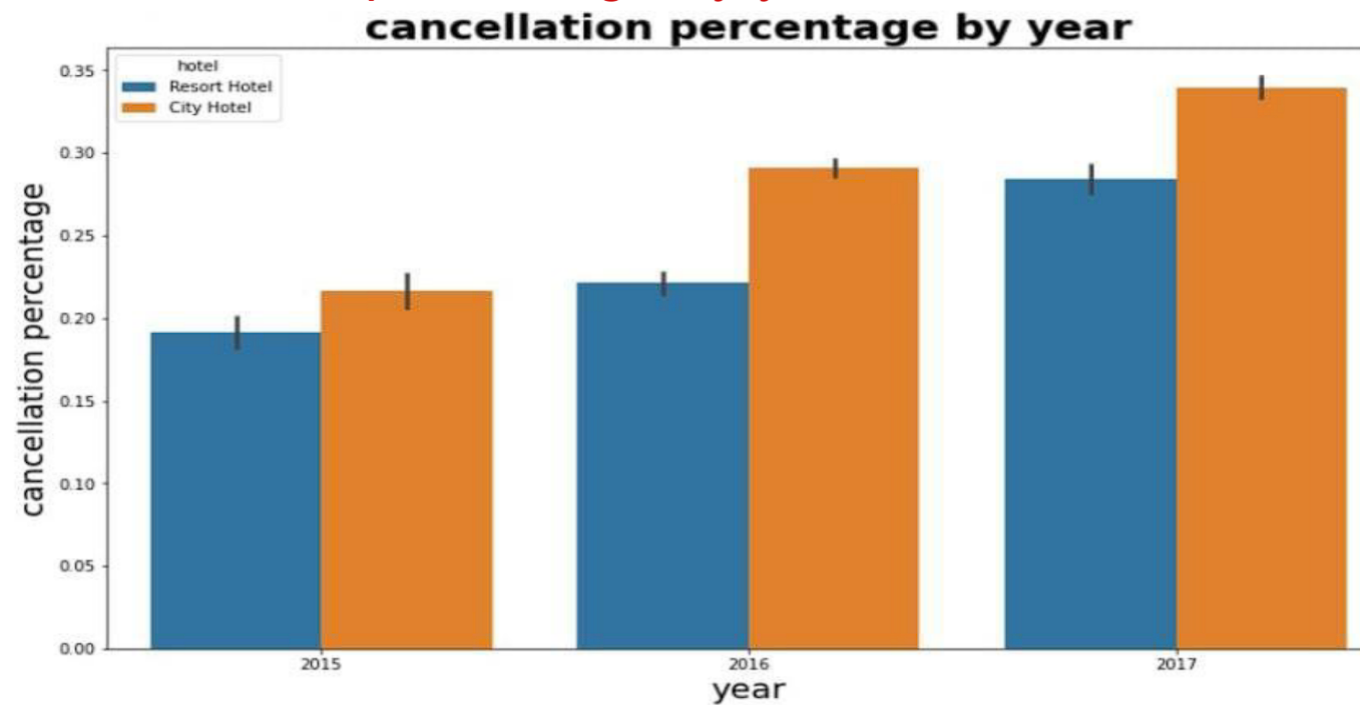
## Number of kids in both hotels



## Heatmap between type of Hotel and arrival month

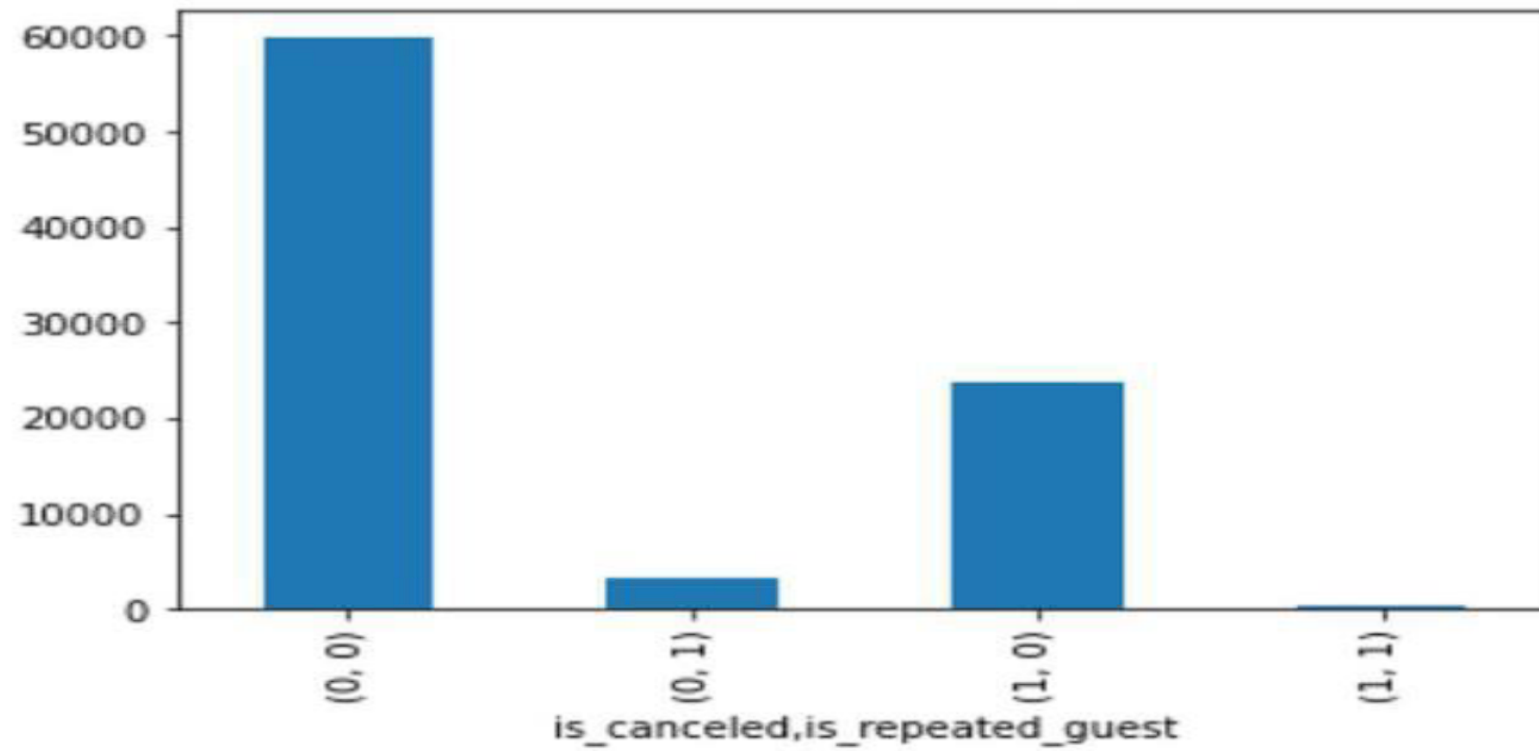


## Cancellation percentage by year



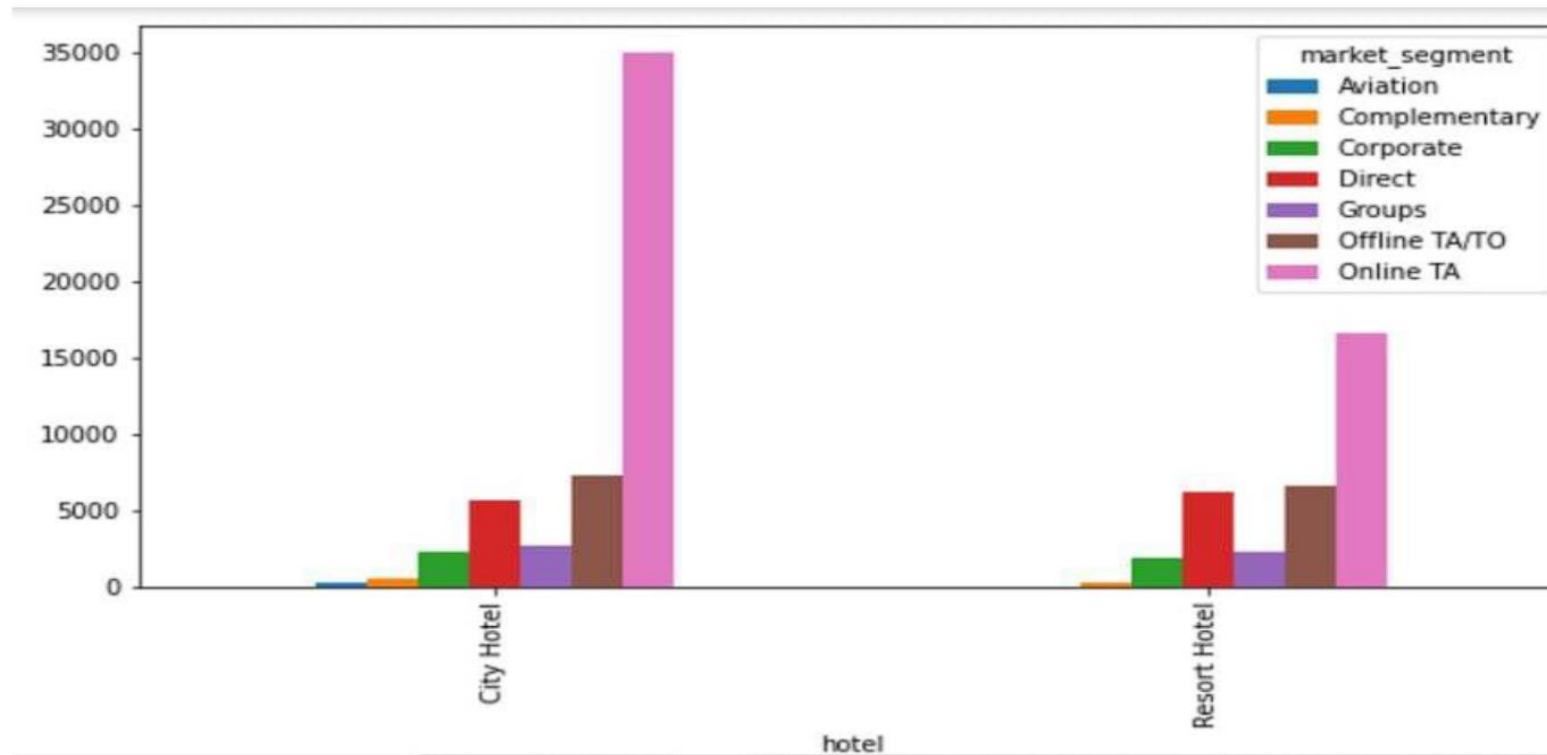


## Plot between cancellation type & repeated guest

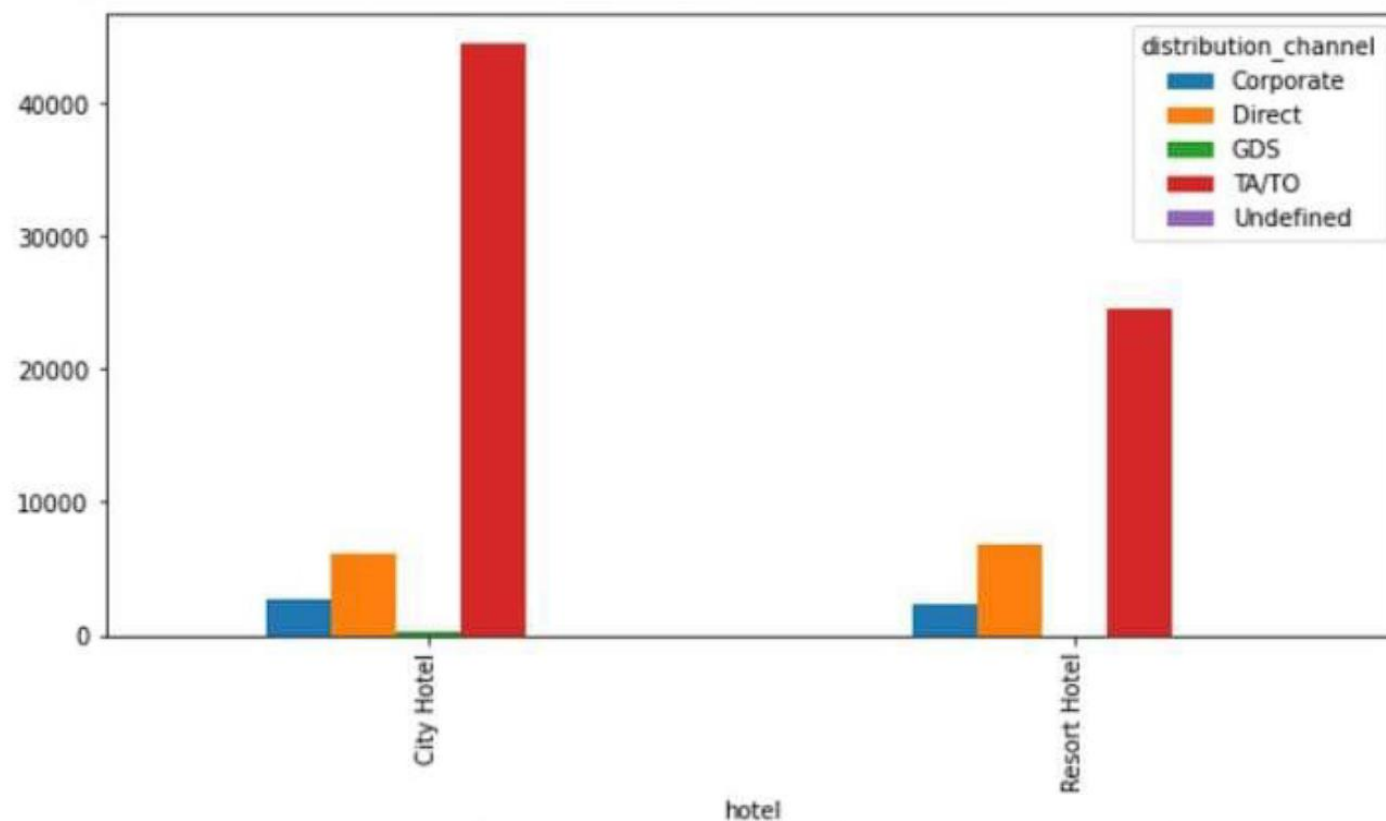


# BIVARIATE AND MULTIVARIATE ANALYSIS

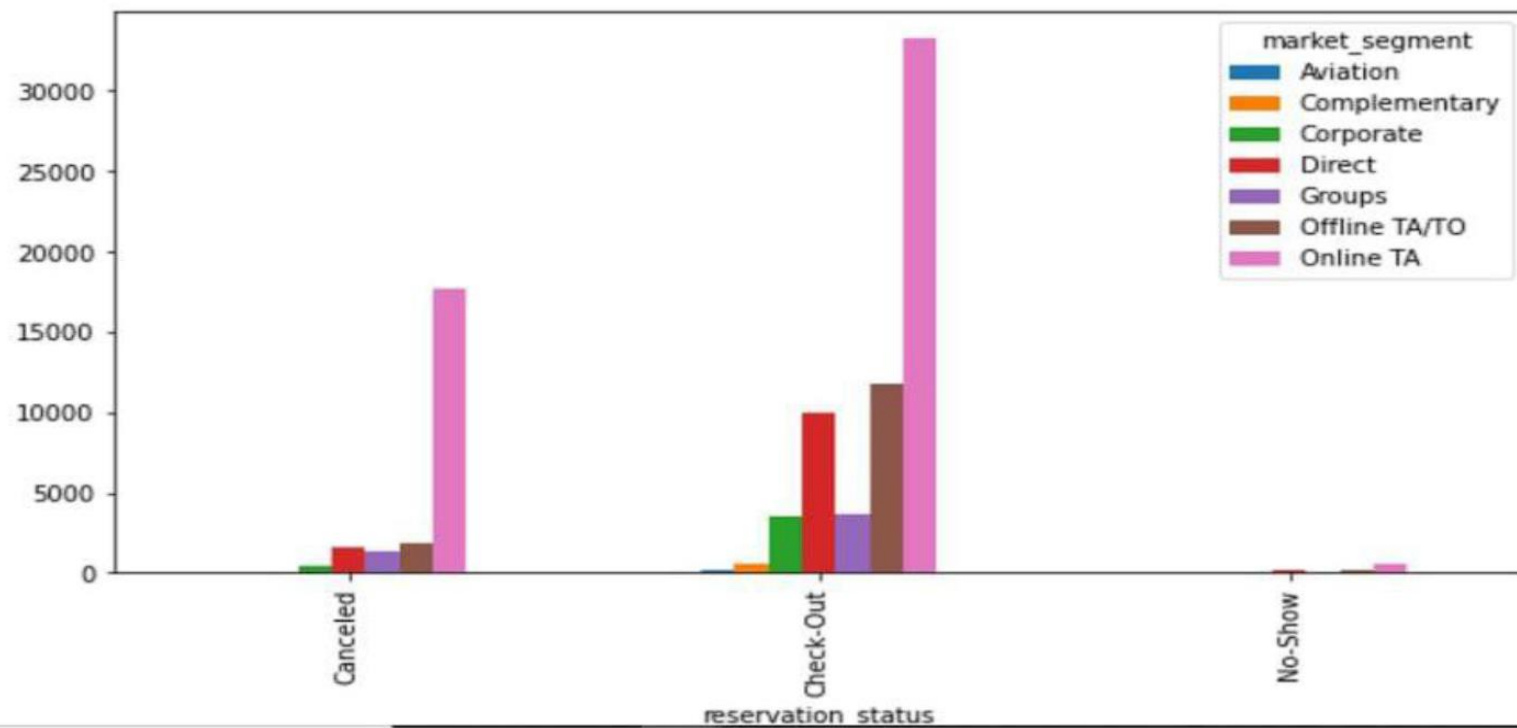
Plot between hotel and market segment



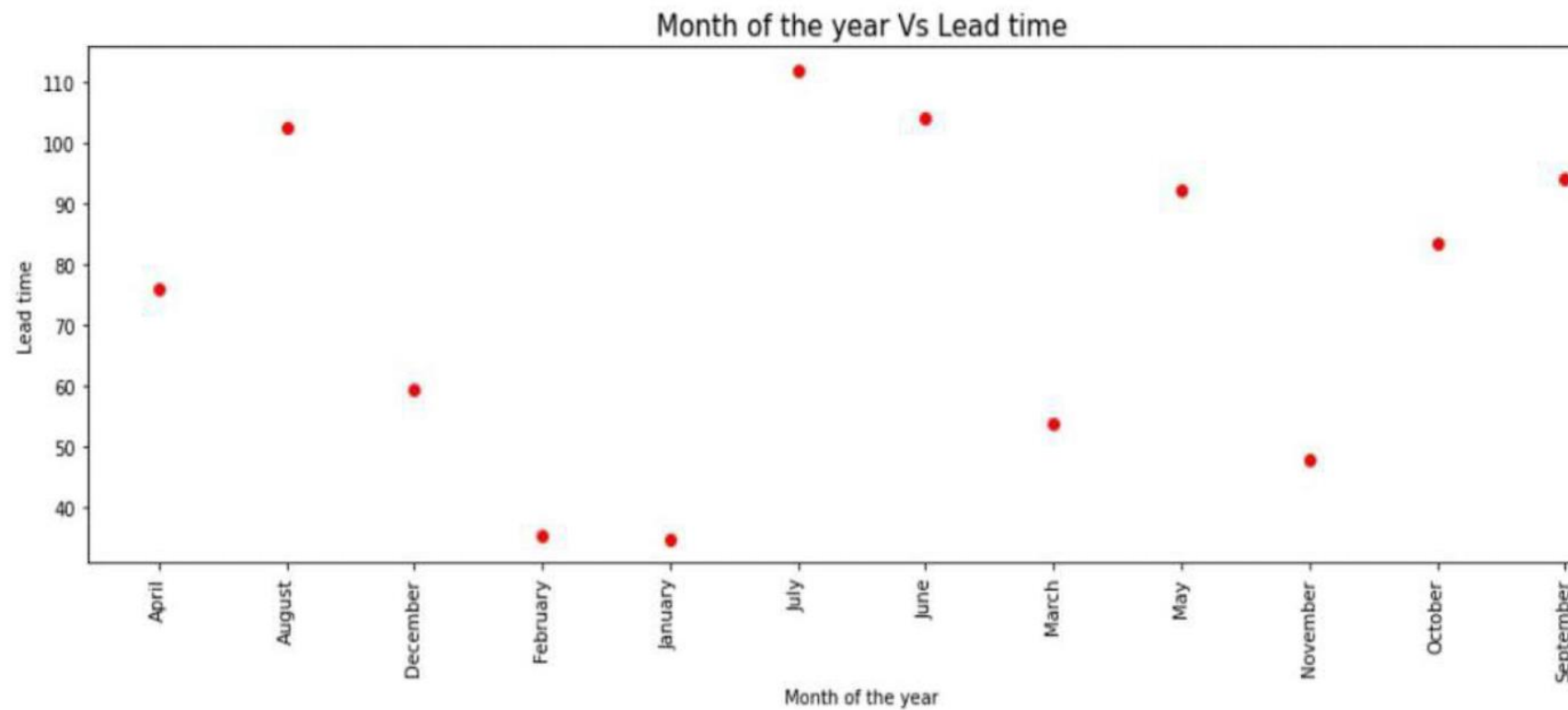
## Plot between hotel and distribution channel



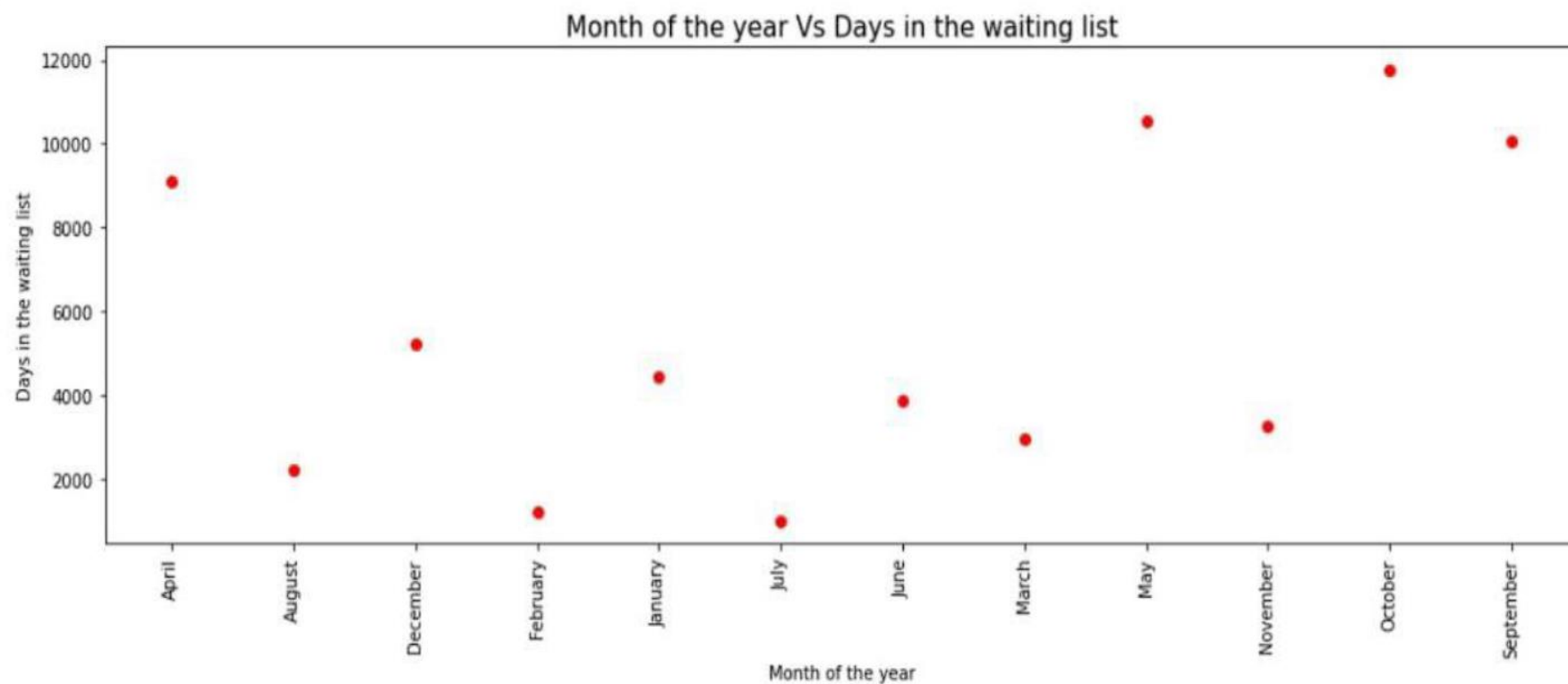
## Plot between resevation status and market segment

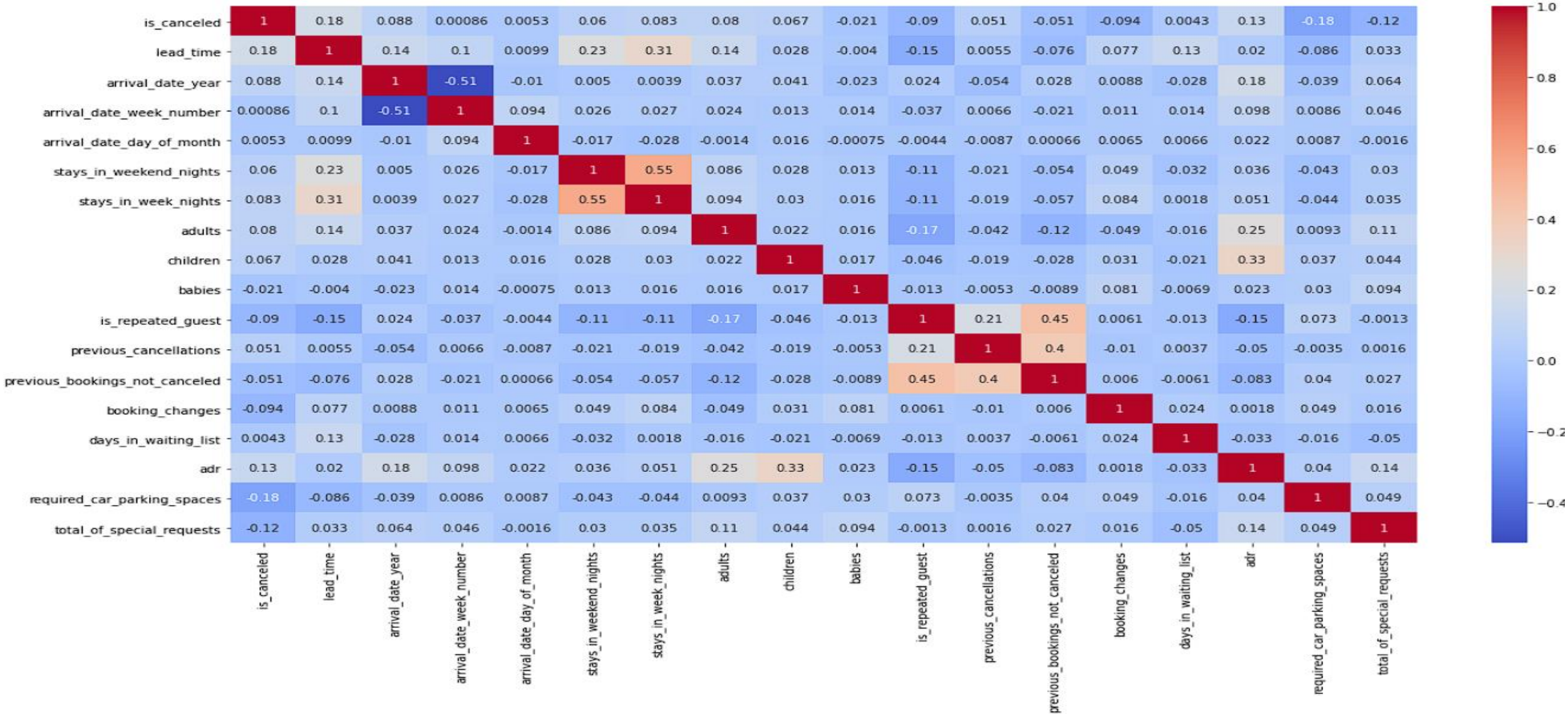


## Plot between Month of year & Lead time



## Plot between Month of year & Days in the waiting list





# CHALLENGES

- The name of the countries was not in the proper format, because of which we are not able to plot the geomap plot.

Company and agent column has lots of duplicate value

- There were many rows with almost similar data
- Lots of null values in the dataset



# CONCLUSION

- Month of August and July receives most no. of booking.
- Booking for city hotels is twice as for resort hotels.
- Repeated customers cancel their hotel in very rare cases.
- Customers coming from aviation industry have very less time i.e. they book urgently
- People with no kids prefer to choose city hotel over resort hotel

# Strategies to counter high cancellations at hotel

- Since we see, our repetitive costumers are most loyal costumers,to maintain them we can provide them with some bonus points,which can be redeem in the next booking
- Month of January and December receives less no. of booking,hotels can offer discounted packages for these months.
- Family with kids prefer resorts , we can provide with holiday family packages.
- Great no. of the bookings are coming from travel agents, so we can provide them some commission

**THANK YOU**